



Co-scheduling HPC workloads on cache-partitioned CMP platforms

Guillaume Aupy, Anne Benoit, Brice Goglin, Loïc Pottier, Yves Robert

► To cite this version:

Guillaume Aupy, Anne Benoit, Brice Goglin, Loïc Pottier, Yves Robert. Co-scheduling HPC workloads on cache-partitioned CMP platforms. [Research Report] RR-9154, Inria. 2018. hal-01719728

HAL Id: hal-01719728

<https://inria.hal.science/hal-01719728>

Submitted on 28 Feb 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Co-scheduling HPC workloads on cache-partitioned CMP platforms

Guillaume Aupy, Anne Benoit, Brice Goglin, Loïc Pottier, Yves Robert

**RESEARCH
REPORT**

N° 9154

February 2018

Project-Team ROMA, TADaaM



Co-scheduling HPC workloads on cache-partitioned CMP platforms

Guillaume Aupy*, Anne Benoit^{†‡}, Brice Goglin*, Loïc Pottier[†],
Yves Robert^{†§}

Project-Team ROMA, TADaaM

Research Report n° 9154 — February 2018 — 33 pages

Abstract: With the recent advent of many-core architectures such as chip multiprocessors (CMP), the number of processing units accessing a global shared memory is constantly increasing. Co-scheduling techniques are used to improve application throughput on such architectures, but sharing resources often generates critical interferences. In this paper, we focus on the interferences in the last level of cache (LLC) and use the *Cache Allocation Technology* (CAT) recently provided by Intel to partition the LLC and give each co-scheduled application their own cache area. We consider m iterative HPC applications running concurrently and answer to the following questions: (i) how to precisely model the behavior of these applications on the cache partitioned platform? and (ii) how many cores and cache fractions should be assigned to each application to maximize the platform efficiency? Here, platform efficiency is defined as maximizing the performance either globally, or as guaranteeing a fixed ratio of iterations per second for each application. Through extensive experiments using CAT, we demonstrate the impact of cache partitioning when multiple HPC application are co-scheduled onto CMP platforms.

Key-words: co-scheduling, HPC application, cache-partitioning, chip multiprocessor (CMP).

* Inria, Labri, Univ. Bordeaux, CNRS, Bordeaux-INP, France

† LIP, École Normale Supérieure de Lyon, CNRS & Inria, France

‡ Georgia Institute of Technology, Atlanta, GA, USA

§ University Tennessee Knoxville, USA

Co-ordonnancement d'applications scientifiques sur des plates-formes multi-coeurs avec partition du cache

Résumé : Ce rapport étudie les techniques de partitionnement de cache pour le co-ordonnancement d'applications scientifiques sur plates-formes multi-coeurs. Nous nous focalisons sur les interférences dans le cache de dernier niveau et utilisons la technologie CAT (*Cache Allocation Technology*) récemment proposée par Intel pour partitionner le LLC et allouer à chaque application sa propre zone de cache. Nous considérons m applications itératives qui s'exécutent simultanément et répondons aux questions suivantes: (i) comment modéliser de façon précise le comportement de ces applications; (ii) combien de coeurs et quelle fraction de cache allouer à chaque application? Notre objectif est de maximiser la performance quand on impose un ratio relatif d'itérations par application, ce qui revient à maximiser le plus petit débit applicatif (pondéré par ces ratios). Ensuite, via un jeu complet d'expérimentations avec CAT, nous montrons l'impact des techniques de partitionnement de cache dans ce contexte, et quantifions le gain qu'on peut en attendre.

Mots-clés : co-scheduling, application scientifique, partitionnement de cache, chip multiprocessor (CMP).

1 Introduction

Co-scheduling applications on a chip multiprocessor (CMP) has received a lot of attention recently [1, 2]. The main motivation is to improve the efficiency of the parallel execution of each application. Consider for instance the Gyoukou ZettaScaler supercomputer, currently ranked #4 in the TOP500 benchmark [3]: it uses PEZY-SC2, a 2048-core processor chip sharing a 40MB last level cache (LLC) [4]: with so many cores at disposal, few applications can efficiently be deployed on the entire computing platform. This is because most application speedup profiles obey Amdahl’s law, which tends to severely limit the number of cores to be used in practice.

The remedy is simple: schedule many applications to execute concurrently; each application receives only a fraction of the total number of cores, hence its parallel efficiency is improved. Which fraction of computing resources should actually be assigned to each application depends on many factors, including speedup profiles, but also external constraints prescribed by the user such as response times or application priorities.

Unfortunately, the remedy comes with complications: when multiple applications run concurrently on a CMP, they compete to access shared resources such as the LLC, and their performance actually degrades. This issue turned out so severe [5, 6] that the name *co-run degradation* has been coined. Modeling and studying cache interferences to prevent co-run degradation has been the object of many studies [7, 8, 9] (see Section 2 on related work for more details).

Intel recently introduced a new hardware feature for cache partitioning called *Cache Allocation Technology* (CAT) [10]. CAT allows the programmer to reserve cache subsections, so that when several applications execute concurrently, each of them has its own cache area. Using CAT, Lo et al. [2] showed experimentally that important gains could be reached by co-scheduling latency-sensitive applications with a strict cache partitioning. In this paper, we also use CAT to partition the LLC into several areas when co-scheduling applications, but with the objective of optimizing the throughput of *in-situ* or *in-transit* analysis for large-scale simulations. Indeed, in such simulations, data is generated at each iteration and periodically analyzed by parallel processes on dedicated nodes, concurrently of the main simulation [11]. If these dedicated nodes belong to the main simulation platform (thereby reducing the number of available cores for simulation), we speak of *in-situ* processing, while if they belong to an auxiliary platform, we speak of *in-transit* processing [12]. In both cases, several applications (various kernels for analysis) have to run concurrently to analyze the data in parallel of the current simulation step. The constraint is to achieve a prescribed throughput for each application, because the outcome of the analysis drives the next steps of the simulation. In the simplest case, each application will have to complete within the time of a simulation step, hence we need to

achieve the same throughput for each application, and maximize that value. In other situations, some applications may be needed only every k simulation steps, with a different value of k per application [13]. This calls for achieving a weighted throughput per application, and for maximizing the minimum value of these weighted throughputs, which dictates the global rate at which the analysis can progress.

The first major contribution of this paper is to introduce a model that characterizes application performance, and to show how to optimally decide how many cores and which cache fraction should be assigned to each application in order to maximize the weighted throughput. The second major contribution is to provide an extensive set of experiments conducted on the Intel Xeon, which assesses the gains achieved by our optimal resource allocation strategy.

The rest of the paper is organized as follows. Section 2 provides an overview of related work. Section 3 details the main framework and all application/platform parameters, as well as the optimization problem. Section 4 presents five co-scheduling strategies, including a dynamic programming approach that provides an optimal resource assignment (according to the model). Section 5 describes the real cache partitioned platform used to perform the experiments. Section 6 assesses the accuracy of the model. Section 7 reports extensive experiments. Finally, Section 8 summarizes our main contributions and discusses directions for future work.

2 Related work

Recent multi-core processors show dozens of cores and a shared cache always larger. In this context, co-scheduling has been extensively studied [1, 2]. The main idea behind co-scheduling is to execute applications concurrently rather than in sequence in order to improve the global throughput of the platform. Indeed, many HPC applications are not perfectly parallel, and it is not beneficial to deploy them on the entire platform: the application speedup becomes too low beyond a given core count. A new trend in large-scale simulations are *in-situ* and *in-transit* approaches, to visualize and analyze the data during the simulation [14]. Basically, the idea behind these approaches is that a new dataset is generated periodically, and we need to run different applications on different parts of this dataset before the next period. In the *in-situ* approach, simulation and analyzes are co-located in the same node, while in the *in-transit* approach, the data analyzes are outsourced onto dedicated nodes [12]. Several studies have shown that large-scale simulations with *in-situ* could benefit from co-scheduling approaches [11, 15]. The difficulty consists in ensuring that the in-situ part processes the data fast enough to avoid slowing down the main simulation, which is directly related to co-scheduling issues: how to partition the resources across the

concurrent analysis applications that share the CMP?

Shared resources include cache, memory, I/O channels and network links, but among potential degradation factors, cache accesses are prominent [16]. Modeling application interferences is challenging, and one way to address this problem is to partition the cache to avoid these potential interferences. Multiple cache partitioning schemes have been designed, through hardware techniques [17, 18, 19] and software techniques [20, 21, 22, 23]. Most of the hardware approaches are efficient with a very low overhead at the execution time, but they suffer from an extra cost in terms of hardware components. In addition, hardware solutions are difficult to implement and often only tested through simulated architectures. On the side of software-based solutions, the most popular is *page coloring*, where physical pages are selected for application allocations so that they end up in specific sections of the cache. Tam et al. [21], showed that important gains can be achieved through a static partitioning of the L2 cache using page coloring. Besides static strategies, dynamic cache partitioning strategies using page coloring have also been studied. In [22], the cache partitioning is refined and adjusted periodically at runtime, with the objective to maximize platform efficiency.

Modeling application interference is a challenging task, Hartstein et al. [24] showed, with the Power Law of cache misses (or the $\sqrt{2}$ rule), how the cache size affects the cache miss ratio. The Power Law states that, if for a baseline cache of size C_0 , the cache miss ratio is equal to m_0 , then for a cache of size C , the cache miss ratio $m = m_0 \left(\frac{C_0}{C}\right)^\alpha$, where α is usually set to 0.5. In a previous work [25] using this law, we were focusing on a static allocation of LLC cache fractions, and core numbers, to concurrent applications as a function of several parameters (cache-miss ratio, access frequency, operation count). We used simulations to assess the performance of our algorithms, because at that time no cache partitioning technologies were available. Intel recently released a new software technique to internally partition the last level cache (LLC), called the *Cache Allocation Technology* (CAT) [10, 2]. In this paper, we use CAT to experiment with a real cache partitioned platform. To the best of our knowledge, this work is the first co-scheduling study for a cache partitioned system (using CAT) with HPC workloads.

3 Model and optimization problem

In this section, we first describe the application model in Section 3.1, and then we formalize the optimization problem in Section 3.2.

3.1 Application model

The objective is to execute m iterative applications A_1, \dots, A_m on P identical cores. The applications are sharing a cache of size C . As explained in

Section 1, new technologies allow us to decide how many cores and which fraction of cache are allocated to each application. Specifically, the cache can be divided into X different fractions. For instance, if $X = 20$, we can give several fractions of 5% of the cache to each application.

Let p_i be the number of cores on which application A_i is executed, and let x_i be the number of fractions of cache assigned to A_i , for $1 \leq i \leq m$. Hence, A_i uses a cache of size $\frac{x_i}{X}C$. We must have $\sum_{i=1}^m p_i = P$ and $\sum_{i=1}^m x_i = X$, i.e., all the cores and the cache fractions are partitioned across the applications.

Given p_i and x_i , an application A_i executes one iteration in time

$$T_i(p_i, x_i) = t_i(p_i) (1 + h_i(x_i)), \quad (1)$$

where $t_i(p_i)$ represents the computation cost and $h_i(x_i)$ the slowdown induced by cache misses in the LLC. Intuitively, the computation cost decreases when p_i increases, and similarly, the slowdown decreases when x_i increases. In this formula, the slowdown incurred by cache misses does not depend on the number of cores assigned to the application. We keep this assumption in our model, and discuss its accuracy in Section 6, where we measure cache misses and refine the model.

Assumption 1 *In the execution time, the slowdown due to cache misses does not depend on the number of cores involved.*

We now detail the model for $t_i(p_i)$ and $h_i(x_i)$.

Computations $t_i(p_i)$. We assume that all applications obey Amdahl's law [26], hence

$$t_i(p_i) = s_i T_i^{seq} + (1 - s_i) \frac{T_i^{seq}}{p_i}, \quad (2)$$

where T_i^{seq} is the sequential time of the application executed with 100% of the cache, and s_i is the sequential fraction of the application.

Cache misses effect $h_i(x_i)$. The most challenging part is to model the slowdown factor $h_i(x_i)$. In chip multiprocessors (CMP), many studies have observed that cache miss ratio follows the Power Law, also called the $\sqrt{2}$ rule [24, 27, 28]. The Power Law of cache misses states that for a cache of size C_{act} , the cache miss ratio r can be expressed as

$$r = r_0 \left(\frac{C_0}{C_{act}} \right)^\alpha, \quad (3)$$

where r_0 represents the cache miss ratio for a baseline cache of size C_0 , and α is a parameter ranging from 0.3 to 0.7, with an average at 0.5. We consider $\alpha = 0.5$ in the following.

We slightly generalize the Power Law formula (with $\alpha = 0.5$) to avoid side effects, and define the slowdown as follows:

$$h_i(x_i) = a_i + \frac{b_i}{\sqrt{x_i}}, \quad (4)$$

where a_i and b_i are constants depending on the application A_i . From Equation (3) with $\alpha = 0.5$, we have $b_i = r_0 \sqrt{\frac{C_0 X}{C}}$ (since $C_{act} = \frac{x_i}{X} C$). In Section 6, we determine a_i and b_i by interpolation, from experimentally measured cache misses, see Table 2.

Finally, when assigning p_i cores and a fraction x_i of the cache, an application A_i executes one iteration in time

$$T_i(p_i, x_i) = t_i(p_i) \left(c_i + \frac{b_i}{\sqrt{x_i}} \right), \quad (5)$$

where $c_i = 1 + a_i$.

3.2 Optimization problem

As stated in Section 1, the goal is to maximize a weighted throughput, since analysis applications may be required at different rates, from every simulation step to every tenth (or more) step [13]. We let β_i denote the weight of application A_i for $1 \leq i \leq m$. Intuitively, β_i represents the number of times that we should execute application A_i at each iteration step. These priority values are not absolute but relative: for $m = 2$ applications, having $\beta_1 = \frac{1}{4}$ and $\beta_2 = 1$ means we execute four times A_2 (at each step) while executing A_1 only once (every fourth step). This is equivalent to having $\beta_1 = 1$ and $\beta_2 = 4$ if we change the granularity of the simulation steps. In fact, what matters is the relative number of executions of each A_i that is required, hence we aim at maximizing the weighted throughput:

- The throughput achieved when executing β_i instances of application A_i is $\frac{1}{\beta_i T_i(p_i, x_i)}$;
- The objective is to partition the shared cache and assign cores such that the total time taken by the slowest application is minimal, i.e., the lowest weighted throughput is maximal.

The weighted throughput allows us to ensure some fairness between applications, and to enforce a better analysis rate of the simulation results whenever the bottleneck is the slowest application. Of course, letting $\beta_i = 1$ lead to maximizing the rate of the analysis when all applications are needed at the same frequency. The optimization problem is formally expressed below:

Definition 1 (CoSched-CachePart) *Given m iterative applications with priorities $(A_1, \beta_1), \dots, (A_m, \beta_m)$ and a platform with P identical cores sharing a memory of size C with X fractions of cache, the COSCHED-CACHEPART*

problem consists in finding a schedule $\{(p_1, x_1), \dots, (p_m, x_m)\}$ such that

$$\text{MAXIMIZE } \min_{1 \leq i \leq m} \left\{ \frac{1}{\beta_i T_i(p_i, x_i)} \right\} \quad \text{SUBJECT TO } \begin{cases} \sum_{i=1}^m p_i = P, \\ \sum_{i=1}^m x_i = X. \end{cases}$$

4 Scheduling strategies

In this section, we introduce several co-scheduling strategies that we will compare via experiments on the Intel Xeon. We start with a (theoretically) optimal schedule, and then present simple heuristics that we use for reference and comparison.

4.1 Optimal solution to CoSched-CachePart

The optimal solution to COSCHED-CACHEPART can be obtained with a dynamic programming algorithm. Let $T(i, q, c)$ be the maximum weighted throughput that can be obtained with applications A_1, \dots, A_i , using q cores and c fractions of cache. The goal is to find $T(m, P, X)$. We compute $T(i, q, c)$ as follows:

$$T(i, q, c) = \begin{cases} \frac{1}{\beta_1 T_1(q, c)} & \text{if } i = 1, \\ \max_{\substack{1 \leq q_i \leq q \\ 1 \leq c_i \leq c}} \left\{ \min \left\{ T(i-1, q-q_i, c-c_i), \frac{1}{\beta_i T_i(q_i, c_i)} \right\} \right\} & \text{otherwise.} \end{cases}$$

We can compute all values in time $O(mPX)$. In practice on the Intel Xeon, $m \leq P = 14$, and $X = 20$, hence the dynamic programming algorithm will execute almost instantaneously. Checking its optimality in practice will assess the accuracy of the performance model. This strategy is called DP-CP (Dynamic Programming with Cache Partitioning).

4.2 Equal-resource assignment

To evaluate the global efficiency of the optimal solution for DP-CP, we compare it to EQ-CP, a simple strategy that allocates the same number of cores and the same number of cache fractions to each application. The algorithm is the following: we start to give $x_i = \lfloor \frac{X}{m} \rfloor$ and $p_i = \lfloor \frac{P}{m} \rfloor$ for all i , then, we give the $P \bmod m$ extra cores one by one to the first $P \bmod m$ applications, and we give the $X \bmod m$ extra cache fractions one by one to the last $X \bmod m$ applications (see Algorithm 1). Doing this, we forbid the case where an application receives an extra core plus an extra fraction of cache, thereby avoiding a totally unbalanced equal assignment.

4.3 Impact of cache allocation

In order to isolate the impact of cache partitioning on performance, we introduce some variants where only the cache allocation is modified:

- DP-EQUAL uses the number of cores returned by the dynamic programming algorithm, hence the same as for DP-CP, but shares the cache equally across applications, as done by EQ-CP.
- We also consider strategies that do not enforce any cache partitioning, but only decide on the number of cores for each application. DP-NoCP uses the same number of cores as DP-CP, and EQ-NoCP uses an equal-resource assignment as in EQ-CP. However, for these two strategies, all applications share the whole cache, i.e., CAT is disabled.

5 Experimental setup

In this section, we first describe the platform and the benchmark applications in Section 5.1. Then in Section 5.2, we explain in details the *Cache Allocation Technology* CAT.

5.1 Platform and applications

The experimental platform is composed of a Dell PowerEdge R730 server with two Intel Xeon E5-2650L v4 processors (*Broadwell* microarchitecture). Each processor contains $P = 14$ cores (with Hyper-Threading disabled) that share a 35MB last-level cache (*Cluster-on-Die* disabled), divided into $X = 20$ slices (or fractions). Nodes run vanilla 4.11.0 kernel with cache partitioning enabled. Only one processor (with 14 cores) is used for the experiments, since the LLC is not shared between processors.

Cache experiments are very sensitive to perturbations, so we take great care to ensure that all experiments are fully reproducible. To avoid perturbations, (i) we average values obtained (like cache misses) over 20 (in Section 6) or 5 (in Section 7) identical runs; (ii) we flush the last-level cache entirely between runs; and (iii) experiments run on a dedicated processor while the code that manages the experiments runs on the other processor. All the data presented in this paper (cache misses, number of floating operations, etc), is obtained with PAPI [29].

Algorithm 1: Equal allocation with cache partitioning

```

1 EQ-CP ( $m, P, X$ ) begin
2   for  $i = 1$  to  $m$  do  $p_i \leftarrow \lfloor \frac{P}{m} \rfloor$ ;  $x_i \leftarrow \lfloor \frac{X}{m} \rfloor$ ;
3   for  $i = 1$  to  $P \bmod m$  do  $p_i \leftarrow p_i + 1$ ;
4   for  $i = 1$  to  $X \bmod m$  do  $x_{m+1-i} \leftarrow x_{m+1-i} + 1$ ;
5 end
```

For validations and performance evaluation, we use six HPC workloads from the NAS benchmarks [30] (see Table 1). We consider only NAS benchmarks from class *A*, as detailed in Table 1.

App	Description
CG	Uses conjugate gradients method to solve a large sparse symmetric positive definite system of linear equations
BT	Solves multiple, independent systems of block tridiagonal equations with a predefined block size
LU	Solves regular sparse upper and lower triangular systems
SP	Solves multiple, independent systems of scalar pentadiagonal equations
MG	Performs a multi-grid solve on a sequence of meshes
FT	Performs discrete 3D fast Fourier Transform

Table 1: Description of the NAS parallel benchmarks.

5.2 Cache Allocation Technology

The Cache Allocation Technology (CAT) [10] is part of a larger set of Intel technologies that are called the Intel Resource Director Technology (RDT) support since the *Haswell* architecture. RDT lets the operating system group applications into classes of service (COS). Each class of service describes the amount of resources that assigned applications can use (see Figure 1). Monitoring of current use of these resources may also be available. Currently, resources can be either an amount of cache or memory bandwidth. In this paper we will only focus on cache resources (CAT), which implements cache partitioning.

The CAT divides the LLC into X slices of cache. Each class of service has a set of slices that applications can use: When reading or writing memory requires to fetch a cache line in the LLC, that cache line must be allocated in the slices available to the class of the current application. However applications may read/modify cache lines that are already available in other slices, for instance when sharing memory between programs in different classes (each cache line can only exist once in the entire cache).

Each slice may only be used by a single class. By default, applications

are placed in the default class (COS_0) which contains slices not used by any other class. The set of slices available to a class is a capacity bit-mask (CBM) of length X . With $X = 20$, if COS_1 has access to the last 4 slices (the top 20% of the LLC), CBM_1 would be set to $0xf0000$.

However, CAT has some technical restrictions:

- Number of slices (CBM length) and classes are architecture dependent (20 and 16 on our platform);
- A CBM cannot be empty (each class of applications must have at least one fraction of cache);
- Bits set in a CBM must be contiguous;
- Slices are not distributed geographically in the LLC. Address hashing ensures spreading of slices over the entire LLC. In other words, $0x10000$ and $0x00001$ CBM should behave exactly the same with respect to locality; there are no NUCA effects (Non Uniform Cache Access).

In this work, we consider a strict cache partitioning, hence each COS contains only one application (and each cache slice is available to a single application).

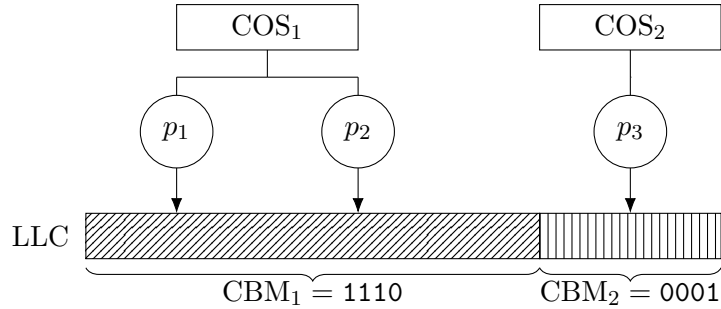


Figure 1: CAT example with 2 classes of service, 3 cores and a 4-bit capacity mask (CBM). First COS has 2 cores and 75% of the LLC, the second class of service has the remaining resources.

6 Accuracy of the model

In this section, we assess the precision of the model developed in Section 3. First, we detail the experimental protocol and explain how to obtain the model parameters for each application in Section 6.1. Then, we study in Section 6.2 the behavior of cache misses on the platform described in Section 5.1, so as to verify whether the Power Law holds for HPC workloads on such architectures. Finally, we study in Section 6.3 the accuracy of the model proposed in Section 3.1 by comparing the expected execution time from Equation (5) to the measured one.

6.1 Experimental protocol

To instantiate the model and check its accuracy, we need to find for each application the value of three parameters used in Equation (5): s_i (sequential fraction), and a_i (or equivalently $c_i = a_i + 1$) and b_i (cache slowdown). To this purpose, we monitor each application with PAPI [29] and use multiple interpolations on the produced data to find the desired constants. To obtain the data on which all the interpolations are based, we proceed as follows: for each application running alone on a dedicated processor, we vary the number of cores from 1 to 14, and for each core number, we also vary the fraction of cache from 5% to 100%. From each run, we collect the number of cache misses and the execution time. The results for all applications are displayed in Table 2.

App_i	a_i	b_i	s_i
BT	-0.0026	0.0287	0.010
CG	-0.0379	0.0474	0
FT	0.0092	0.0129	0.016
LU	-0.0247	0.0275	0.020
MG	0.0460	0.0073	0.065
SP	-0.0110	0.0254	0.018

Table 2: s_i , a_i and b_i obtained by interpolation from the data produced by measurements (averaged on the core numbers, according to Assumption 1).

6.2 Accuracy of the Power Law

Figure 2 shows the evolution of cache miss ratios for the six applications depending on the number of cores and cache fraction. We observe that for most applications, the cache miss ratio increases with the number of cores for small cache fractions, while it does not vary significantly with the number of cores for higher cache fractions. Therefore, these results verify the assumption about the relation between number of cores and cache misses (Assumption 1).

On Figure 3, we study the evolution of cache miss ratios for each considered application, running alone with a single core. We do not look at cache fractions below $x = 3$ (or 15%) because, according to our experiments, it shows irrelevant results due to cache contention. We observe that the Power Law with $\alpha = 0.5$ suits well the behavior of compute-intensive benchmarks CG, BT, LU and SP, but struggles to model memory/communication-intensive applications like MG and FT.

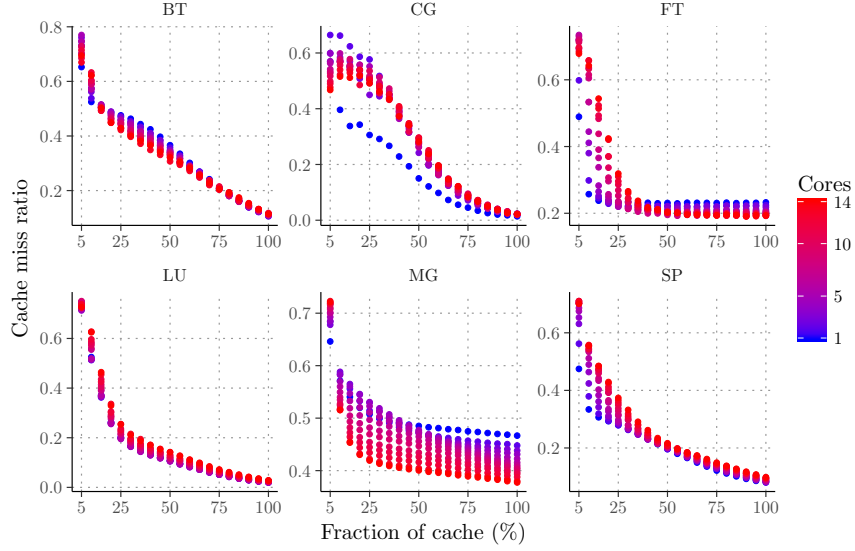


Figure 2: Evolution of cache miss ratio when the cache fraction x_i is ranging from 1 to 20 (i.e., from 5% to 100%) and the number of cores p_i is ranging from 1 (blue) to 14 (red).

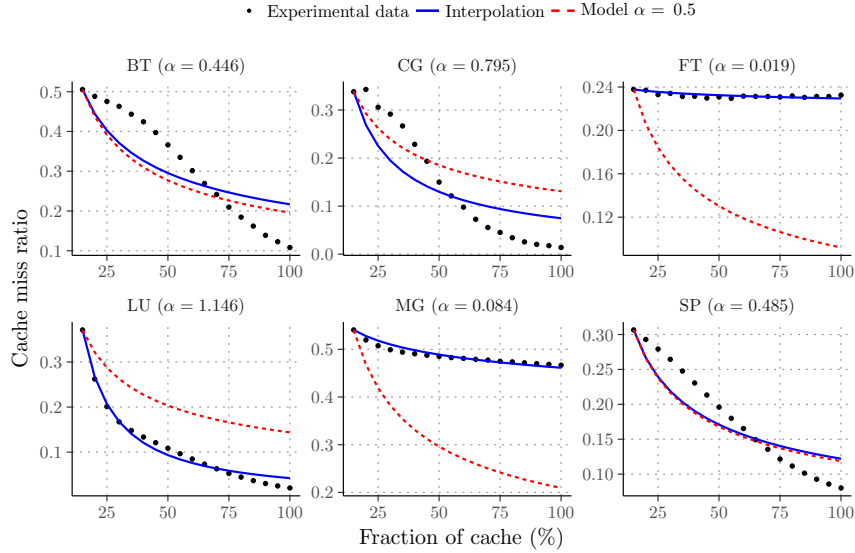


Figure 3: Comparison between the predicted cache miss ratio given by the Power Law with $\alpha = 0.5$ in red, the best found α parameter in blue and the measured cache miss ratio in black. Applications run alone on the platform with 1 core.

6.3 Accuracy of the execution time

We aim at verifying the accuracy of the execution time predicted by the model. Figure 4 shows, for each application, the comparison between the measured execution time and the model, when the application runs alone on the platform (no co-scheduling here). In Figure 4, the number of cores varies from 1 to 14 while the cache fraction is fixed at $x = 3$ (or 15%).

Figure 5 shows the relative error between predictions and the real data. The relative error is defined as

$$E_i(p_i, x_i) = \frac{|T_i(p_i, x_i) - T_i^{real}(p_i, x_i)|}{T_i^{real}(p_i, x_i)},$$

where $T_i^{real}(p_i, x_i)$ is the measured execution time on the cache partitioned platform for application A_i with p_i cores and x_i fractions of cache. We observe that our model predicts execution times rather well for LU, BT, CG and MG, with less than 25% of error for worst cases. For FT, the model is accurate for $x_i \geq 6$ (30%) and $p_i \leq 10$, with a relative error below 15%, but the model loses accuracy for small cache fractions and high number of cores. For SP, we have the same observation, the model is not accurate for a number of cores larger than 8 if the cache fraction is below 50% (the red part in the Figure 5). This is due to a specific behavior of FT and SP: their execution times tend to become constant after a certain core threshold (see Figure 4), while the model expects a strictly decreasing execution time. For both applications, this constant plateau is not due to Amdahl's law (both FT and SP are enough parallel to scale up to 14 cores), hence a contention effect (either from the cache or the memory bandwidth) is probably behind this constant level in performance. Another reason to explain these mispredictions when the number of cores increases, is Assumption 1, which states that the number of cores does not impact LLC cache misses, which is not true for all applications in practice.

7 Results

To assess the performance of the scheduling strategies of Section 4 and to evaluate the impact of cache partitioning on co-scheduling performance, we conduct an extensive campaign of experiments using a real cache partitioned system.

7.1 Experimental protocol

The platform and the applications used for all the experiments are described in Section 5. Recall that we consider iterative applications, hence we have modified their main loop such that each of them computes for a duration T . We choose a value for T large enough to ensure that each application

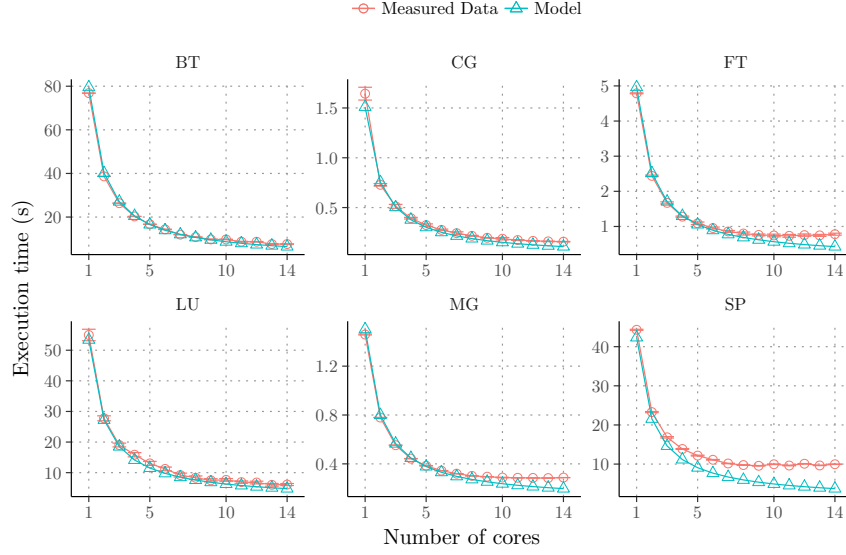


Figure 4: Comparison between predicted execution time by the model and measured execution time, when varying the number of cores up to 14 and with a cache fraction set to 15%.

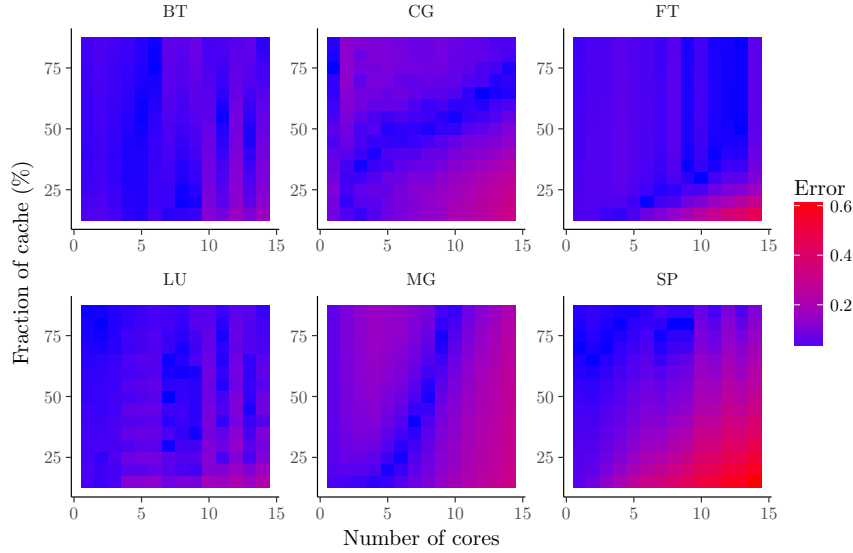


Figure 5: Heat-map of the relative error between the model predictions and the measured execution times when the cache fraction is varying from 15% to 85% and the number of cores from 1 to 14.

reaches the steady state with enough iterations (for instance, $T = 3$ minutes for small applications like CG, FT, MG and $T = 10$ minutes for the others). If a co-schedule contains both small and big applications, we use $T = 10$ minutes for all applications. In addition, for all the following experiments, we use 12 cores out of the 14 available, to avoid rounding effects when we co-schedule a number of applications that is not divisible by the number of cores.

Evaluation framework: To study the performance of the different algorithms under the objective COSCHED-CACHEPART, we measure the time for one iteration of A_i : $T_i = \frac{T}{\#iter_i}$, where $\#iter_i$ is the number of iterations of application A_i during T . Then, we measure $\min_i \frac{1}{\beta_i T_i}$.

To understand the performance of the different algorithms, we are interested by the relative speeds of each application with respect to the others. Intuitively for all i, j , we would like $\beta_i T_i = \beta_j T_j$. Hence we further study the relative error:

$$E_r = \sum_{i \neq j} \left| \frac{\beta_i T_i}{\beta_j T_j} - 1 \right|. \quad (6)$$

7.2 Impact of cache partitioning

The first step is to assess the impact of cache partitioning (CP) on performance. To this purpose, we co-schedule two applications, so we have three combinations (CG+MG, CG+FT, FT+MG). For all i, j , we set the number of cores for A_i and A_j to six, and we vary the fraction of cache allocated to A_i from 5% to 95% while, at the same time, the cache fraction of A_j is varying from 95% to 5%. The y -axis represents the aggregated number of iterations executed by all applications. We run the applications both with CP enabled, and CP not enabled. Figure 6 shows the impact of CP for CG+MG: we can see that when CG has more than 35% of the cache, CP outperforms the version without CP. The impact of CP lies in the behavior of each application, more specifically their data access pattern. CG is a compute intensive application with an irregular memory access pattern, while MG is a memory intensive application. More specifically, MG does not take a great benefit for more cache after 35%, while the performance of CG greatly depends on the cache size (for more details on application behaviors, see Figure 2). Without a cache partitioning scheme, by reading/writing a lot of different cache lines, MG will often evict CG cache lines, resulting into a performance degradation of both applications.

Figure 7 shows the impact of CP for CG+FT. In this case, we note a small improvement when CG has 80% of the cache. The reason behind this improvement is that FT is more communication intensive (all-to-all communication) than strictly memory intensive, hence the gain obtained by CP is less important than for CG+MG. Since we consider only one processor,

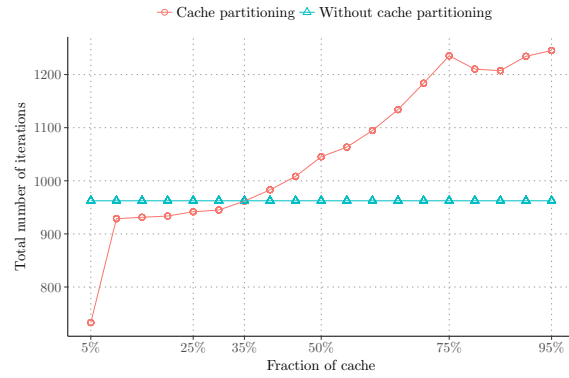


Figure 6: CG and MG with 6 cores each, CG has 5% of the cache while MG has the remaining 95%, then CG has 10% and MG 90% and so forth.

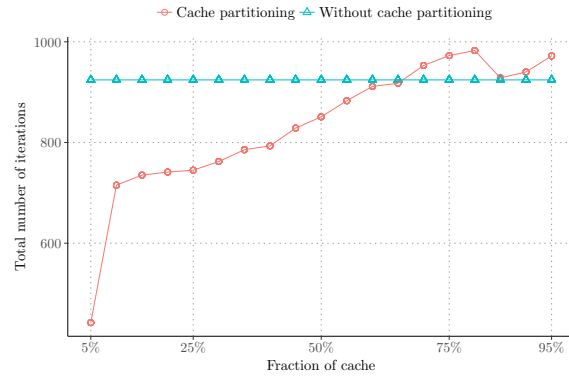


Figure 7: CG and FT with 6 cores each.

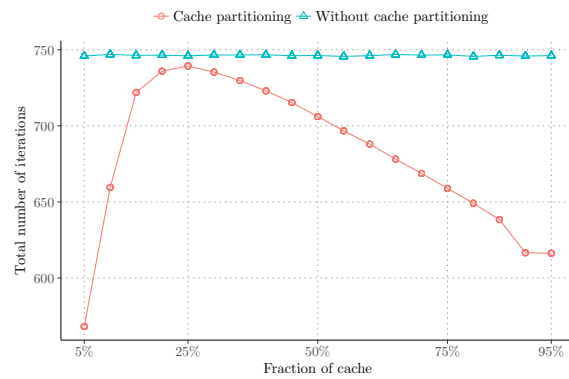


Figure 8: FT and MG with 6 cores each.

the applications that run are the shared memory version (OpenMP), and in that context, the impact of cache on communications is small.

Finally, Figure 8 presents the result for the last combination FT+MG. The cache partitioning is not efficient for that combination of two memory and communication intensive applications. If FT has 25% and MG has 75%, then CP can almost achieve the same performance as without CP. This inefficiency is mostly due to the memory intensive and communication intensive behaviors of both applications involved, none of them needs a strict cache partitioning, since their use of the cache varies during iterations.

Summary: The cache partitioning is very interesting when compute-intensive and memory-intensive application are co-scheduled (important gain, up to 25%, for CG+MG, small gain for CG+FT). On the contrary, FT and MG together perform badly with the cache partitioning enabled, these applications do not benefit from the cache to improve their execution time by iteration. Hence, the behavior of applications has a strong impact on the global performance of cache partitioning, and in general, co-scheduling applications with the same behavior results in degraded global performance when using CP.

7.3 Co-scheduling results with two applications

Now that we have demonstrated the interest of cache partitioning, we study the performance of the scheduling strategies of Section 4. Recall that the COSCHED-CACHEPART optimization problem aims at maximizing the minimum weighted throughput among co-scheduled applications. Considering two applications (A_i, A_j) , for β_i iterations of A_i , we aim at performing β_j iterations of A_j . To avoid some cache effects that appear when the cache area is too small, we set the minimum cache fraction allocated to each application to three (each application has at least 15% of the cache), while the minimum number of cores per application is set to one. We use three different ways to present the result for each studied combination: (i) the objective we want to maximize (minimum weighted throughput), (ii) the ratio of iterations done, and (iii) the relative error defined in Equation (6).

CG+MG: On Figure 9, we see what is the minimum throughput achieved by each method for CG+MG. The weight β associated to MG varies from 0.25 to 4. The algorithms based on dynamic programming DP-CP, DP-EQUAL and DP-NoCP outperform both equal-resource assignment heuristics EQ-CP and EQ-NoCP. In this scenario, the cache partitioning provides a good performance improvement, since on average DP-CP outperforms DP-NoCP.

Figure 10 shows the ratio of iterations for CG+MG. Ideally, we would like to obtain $\beta_{CG}T_{CG} = \beta_{MG}T_{MG}$, the dashed black line represents that opti-

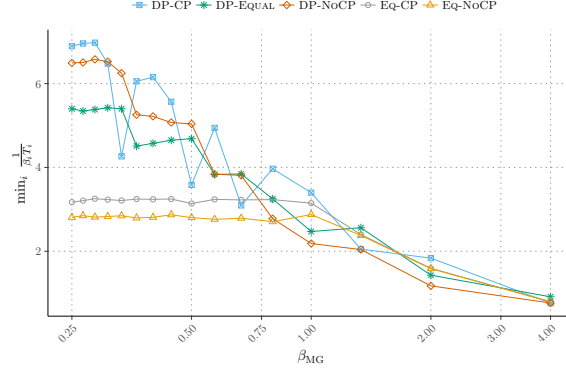


Figure 9: Minimum throughput obtained with CG and MG when β_{MG} is varying from 0.25 to 4 (higher is better).

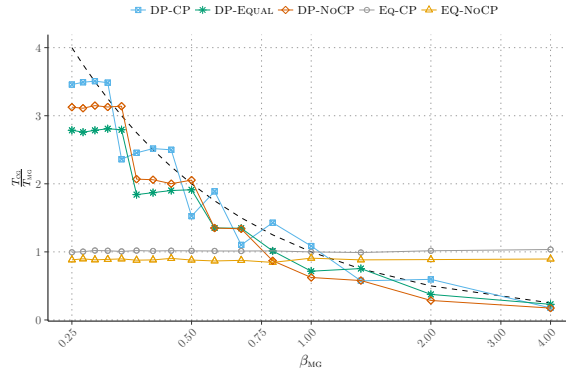


Figure 10: Iteration ratio done by CG and MG when β_{MG} is varying from 0.25 to 4 (closer to the dashed black line is better).

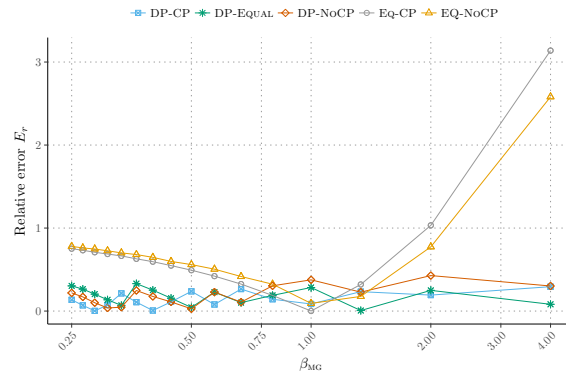


Figure 11: Relative error from the objective for CG and MG when β_{MG} is varying from 0.25 to 4 (lower is better).

mal iteration ratio. First, note that EQ-CP and EQ-NoCP show constant results because they do not depend on weight, but EQ-CP performs better (even without a clever algorithm, cache partitioning helps). Second, we observe that DP-CP is the closest (on average) to the ideal line, hence the cache partitioning really helps here.

Finally, Figure 11 presents the relative error, as defined in Equation (6). We observe that DP-CP, DP-NoCP and DP-EQUAL exhibit the same relative error, near to zero, while EQ-CP and EQ-NoCP present an important error.

CG+FT: In Figure 12, we observe that DP-CP, DP-EQUAL and DP-NoCP outperform EQ-CP and EQ-NoCP when β_{FT} is larger than 0.5. Only, DP-NoCP outperforms EQ-NoCP all the time. When β_{FT} is smaller than 0.5, the two variants without cache partitioning perform better than the two versions with cache partitioning. As explained in Section 7.2, due to its communication-intensive behavior, FT will not benefit a lot from cache partitioning techniques. Figure 13 presents the iteration ratio (i.e., the fairness among co-scheduled applications) when we co-schedule CG+FT: DP-CP, DP-EQUAL and DP-NoCP exhibit good performance, and we are very close to the dashed line that represents the ideal iteration ratio to reach. On Figure 14, we observe the relative error: EQ-CP and EQ-NoCP show an important relative error as expected, and DP-CP, DP-EQUAL and DP-NoCP show the same good performance, very close to zero.

MG+FT: Figure 15 presents the results obtained for MG+FT. DP-CP, DP-EQUAL and DP-NoCP outperform EQ-CP and EQ-NoCP, except for β_{FT} lower than 0.50. For both DP-CP and EQ-CP, the cache partitioning does not bring a important improvement. The main reason is that co-scheduling one memory and one communication intensive application is not very efficient (see Section 7.2). Figure 16 shows that DP-CP, DP-EQUAL and DP-NoCP perform well, very close to the ideal iteration ratio (the dashed line). On Figure 17, we note that the relative error is close to zero for DP-CP, DP-EQUAL and DP-NoCP, while (logically) the relative error is larger for EQ-CP and EQ-NoCP.

LU, SP, BT co-scheduled with MG: Figure 18, Figure 19 and the Figure 20 show the minimum throughput (on the left) and the error norm (on the right) obtained by co-scheduling, respectively, BT+MG, LU+MG and SP+MG. For the minimal throughput (on the left of each figure), both results are quite similar, all variants based on our algorithm DP-CP outperform EQ-CP and EQ-NoCP. The cache partitioning does not bring a significant gain in this scenario, but DP-CP is always better than DP-NoCP. We observe that DP-EQUAL is always worst than DP-CP and DP-NoCP,

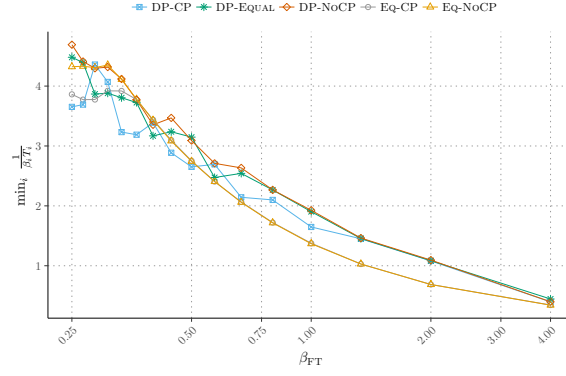


Figure 12: Minimum throughput obtained with CG and FT when β_{FT} is varying from 0.25 to 4.

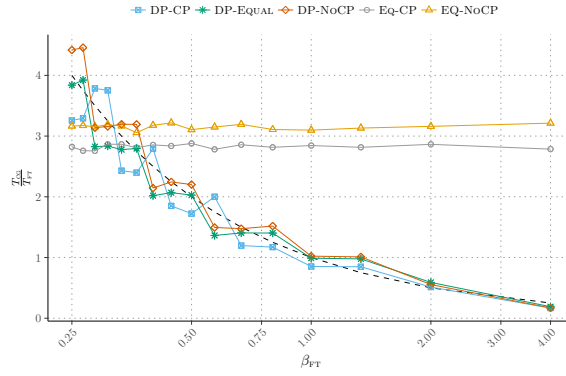


Figure 13: Iteration ratio done by CG and FT when β_{FT} is varying from 0.25 to 4.

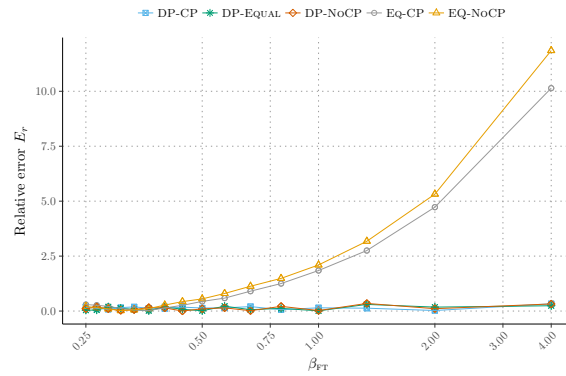


Figure 14: Relative error from the objective for CG and FT when β_{FT} is varying from 0.25 to 4.

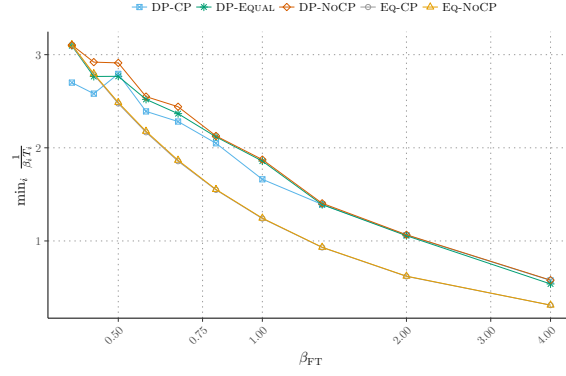


Figure 15: Minimum throughput obtained with MG and FT when β_{FT} is varying from 0.25 to 4.

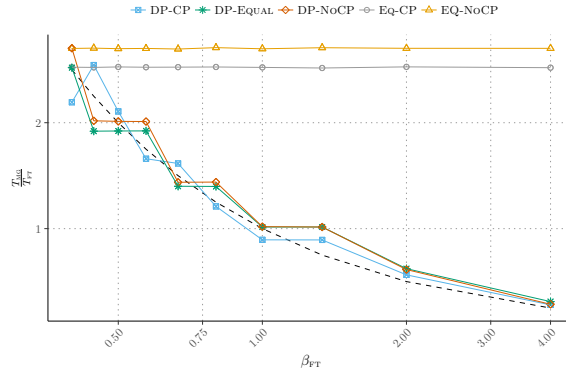


Figure 16: Iteration ratio done by MG and FT when β_{FT} is varying from 0.25 to 4.

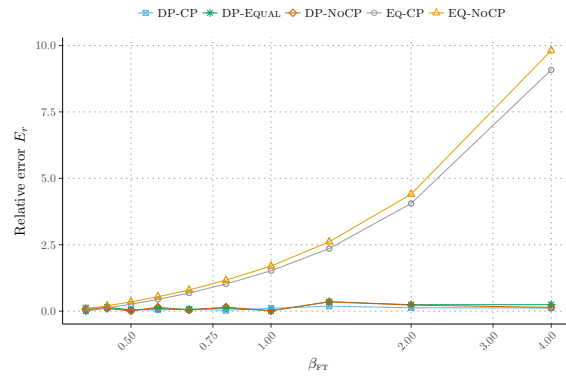


Figure 17: Relative error for MG and FT when β_{FT} is varying from 0.25 to 4.

which means that doing a naive cache partitioning (an equal one in that case) can lead to important performance degradation. For this scenario, because of the high values of the relative error (respectively 0.25 and 0.4 for the best cases), we only present the relative error norm. Indeed, BT, LU and SP are much bigger than MG in terms of number of operations (by roughly 10^3), hence it is impossible to do, for example, four times more iterations of MG than iterations of LU without a very large value of T .

Special case: CG and MG when each application has six cores:

We are now interested into a special case: how the cache will affect co-scheduling performance. All applications have the same number of cores (six in our case), so only the cache is available to increase performance. Figure 21 shows the global performance of all methods. Obviously, only DP-CP takes advantage of this scenario because only this method can choose how to partition the cache. If β_{MG} is smaller than 1, it means that we have to compute more CG than MG, and in that case, the cache has a strong effect (up to 25% improvement with cache partitioning enabled). With this scenario, we are able to isolate which part of performance relies on cache effect. Figure 22 depicts the iteration ratio achieved with an equal number of cores for each application. We observe that with only the cache, it is hard to enforce the required ratio of the number of iterations, according to the values of the β_i . Figure 23 represents the relative error between the ideal iteration ratio and the iteration ratio obtained with each method. Note that the relative error is high for every method, but the error of DP-CP is the smallest.

Summary: The model is accurate enough to enforce that the corresponding optimal DP algorithm performs well: in most cases, DP-CP, DP-EQUAL and DP-NoCP outperform EQ-CP and EQ-NoCP. On the cache partitioning side, when co-scheduling CG and MG, the cache partitioning is really interesting to isolate applications that pollute the cache, such as MG. Figure 21 clearly shows the impact of cache on performance when the number of cores is set for each application. In the worst cases, for instance with FT and MG, the cache partitioning does not improve performance, but does not degrade it either.

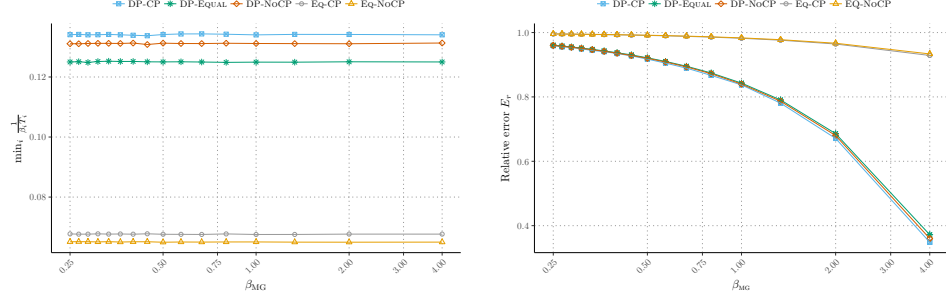


Figure 18: Minimum throughput and relative error for BT+MG.

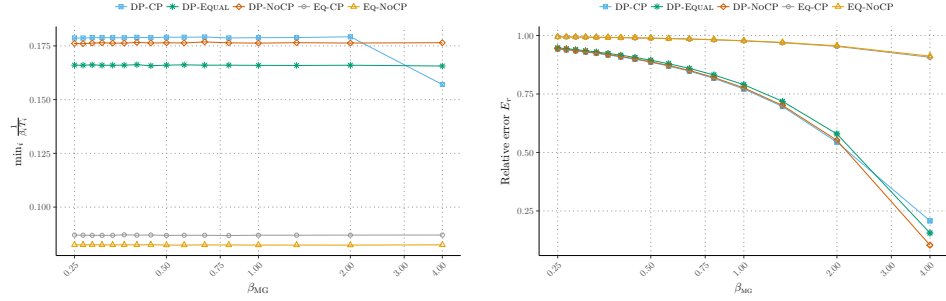


Figure 19: Minimum throughput and relative error for LU+MG.

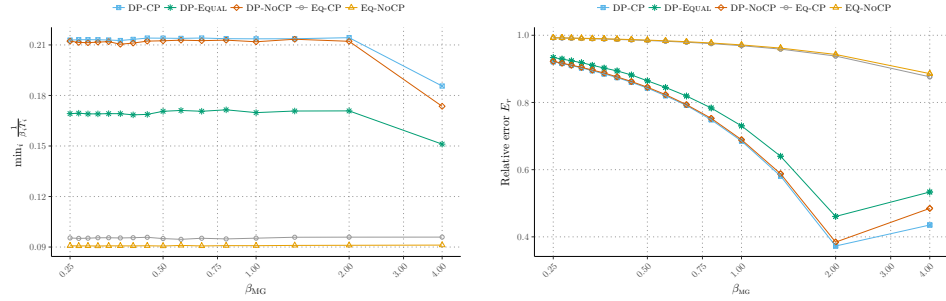
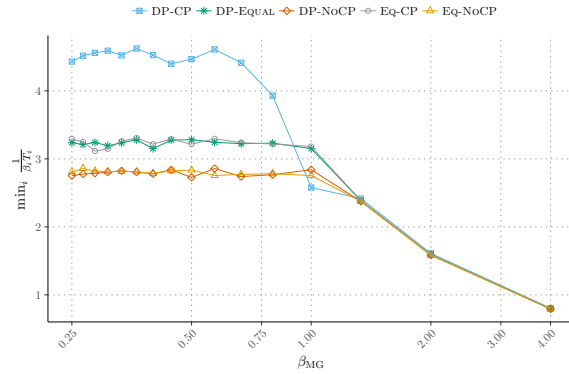


Figure 20: Minimum throughput and relative error for SP+MG.

Figure 21: Minimal throughput obtained with CG and MG when β_{MG} is varying from 0.25 to 4 when both applications have six cores.

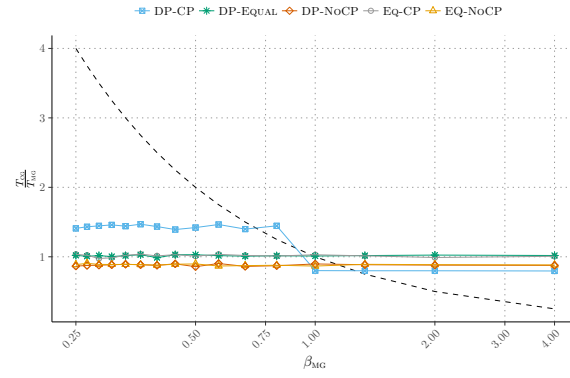


Figure 22: Iteration ratio done by CG and MG when β_{MG} is varying from 0.25 to 4 when both applications have six cores.

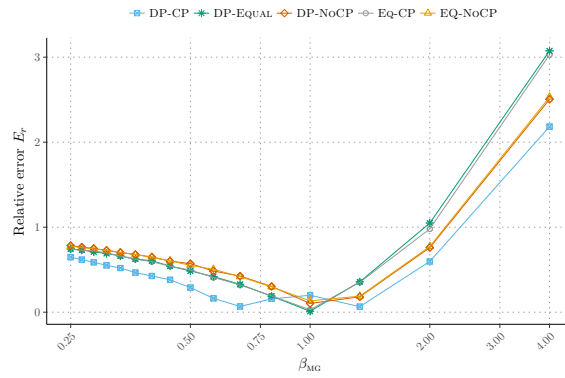


Figure 23: Relative error for CG and MG when β_{MG} is varying from 0.25 to 4 when both applications have six cores.

7.4 Co-scheduling results with three applications

In this section, we present the results with three co-scheduled applications. Similarly to the case with two applications, with three applications (A_1, A_2, A_3) , only β_3 is ranging from 0.25 to 4, while $\beta_1 = \beta_2 = 1$. First, we focus only on co-schedules with CG and MG, because they are very interesting applications to study. Second, we study all combinations of co-scheduling with CG, FT and MG. We do not look at the iteration ratio in this section, but focus on minimum throughput and relative error.

2CG+MG: Figure 24 shows the minimum throughput obtained when we co-schedule 2CG+MG, while the weight associated to MG is ranging from 0.25 to 4. Note that it is interesting to co-schedule multiple copies of the same application (two CGs in this scenario) in order to improve the global efficiency, when this application exhibits a speedup profile with limited gain from adding extra cores and/or extra fractions of caches. We observe that the scheduling strategies building on the dynamic programming algorithm DP-CP, DP-EQUAL and DP-NoCP outperform EQ-CP and EQ-NoCP. In addition, cache partitioning shows a great interest here: DP-CP exhibits a gain around 15% on average over DP-NoCP and DP-EQUAL. The relative error is also depicted on the right. Recall that ideally, we would like to have $\beta_i T_i = \beta_j T_j$ for all i, j (see Equation (6)). We observe that the method that is the closest to zero is DP-CP, confirming the strong interest of cache partitioning.

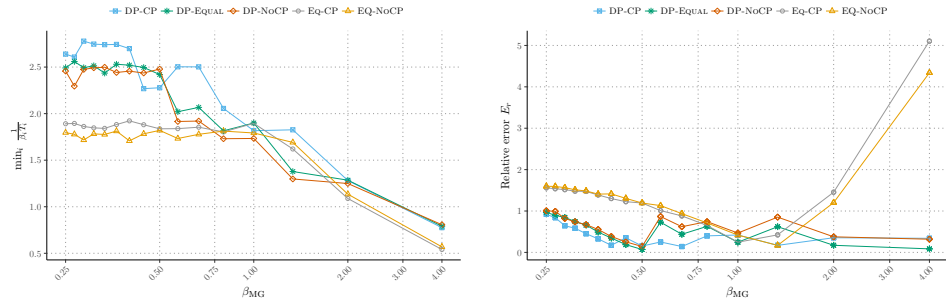


Figure 24: Minimum throughput and relative error for 2CG+MG.

2MG+{CG, BT, LU, SP}: Figure 25 presents the minimal throughput obtained by each method when we co-schedule 2MG+CG, where the weight of CG is ranging from 0.25 to 4. Again, the DP-based strategies DP-CP, DP-EQUAL and DP-NoCP exhibit good performance for β_{CG} smaller than 0.50, but they suffer from a lack of performance when β_{CG} is between 0.50 and 1. When β_{CG} is larger than 1, DP-CP becomes the best method again. On the right of Figure 25, we can see the confirmation that the proposed dynamic programming algorithm is the method that minimizes the best the

relative error, even though the cache partitioning with DP-CP and DP-EQUAL does not bring any clear advantage in this scenario. This is mainly due to the fact that the application with the varying weight is a compute-intensive application, co-scheduled with two memory-intensive applications. According to our experiments, when compute-intensive applications are outnumbered by memory-intensive applications, the cache partitioning is often less efficient.

Figure 26, the Figure 27 and the Figure 28 also presents, the minimal throughput obtained when we co-schedule, respectively, 2MG+BT, 2MG+LU and 2MG+SP. 2MG co-scheduled with BT, LU or SP lead to the same behavior for the minimum throughput and the relative error, the variants based on our dynamic algorithm DP-CP, DP-EQUAL and DP-NoCP perform better than EQ-CP and EQ-NoCP. The error norm, for the three cases, is very important. The reason behind the important values of the error norm is that MG is very small compared to LU, BT and SP.

CG+MG+FT: Figure 29 shows the minimum throughput obtained when co-scheduling the three different applications, while varying only the weight β_{FT} of FT. We observe that the performance of the three DP-based algorithms is close to the performance obtained with the equal-resource assignment for β_{FT} smaller than 0.5, but for the other cases, DP-CP and all its variants outperform EQ-CP and EQ-NoCP. The relative error leads to the same conclusion: DP-CP, DP-NoCP and DP-EQUAL are much closer to zero than EQ-CP and EQ-NoCP, especially when β_{FT} is larger than 0.5.

Next, Figure 30 is the counterpart of Figure 29 when varying only the weight β_{MG} of MG. The results obtained by the DP-based algorithms are very good with an average gain around 50% over the EQ-CP variants, especially when β_{MG} is below 1. We note that the cache partitioning does not take advantage of this scenario, DP-CP shows degraded performance compared to DP-NoCP. For the relative error, the method that performs best is DP-CP, close to DP-NoCP and DP-EQUAL though.

Finally, Figure 31 is the counterpart of Figures 29 and Figure 30 when varying only β_{CG} . The behavior of all DP-CP variants is interesting: for $0.25 \leq \beta_{CG} \leq 0.44$, the resource allocation, both for cores and cache, does not change, resulting into the decreasing of the minimum weighted throughput when β_{CG} is increasing (so $\frac{1}{\beta_{CG}T_{CG}}$, which is actually the minimum here, is decreasing). At $\beta_{CG} = 0.5$, the allocation of resources changes for DP-CP variants (more and more resources are allocated to CG, in order to fit the increasing requirement). We observe that DP-CP, DP-EQUAL and DP-NoCP logically outperform EQ-CP and EQ-NoCP to maximize the minimum weighted throughput among the co-scheduled applications. However, the cache partitioning does not help in this scenario, mainly because we vary the weight of the only compute-intensive application. In terms of rel-

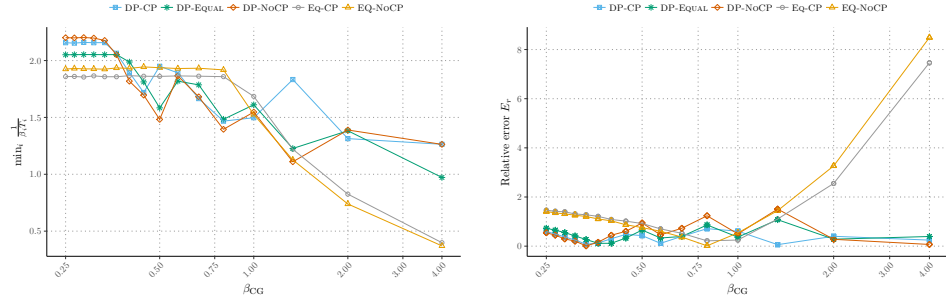


Figure 25: Minimum throughput and relative error for 2MG+CG.

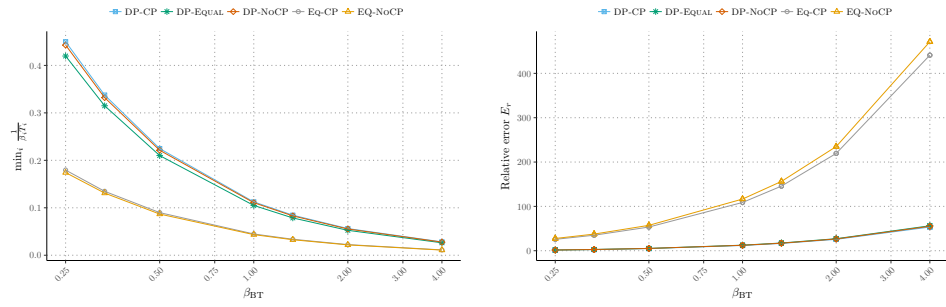


Figure 26: Minimum throughput and relative error for 2MG+BT.

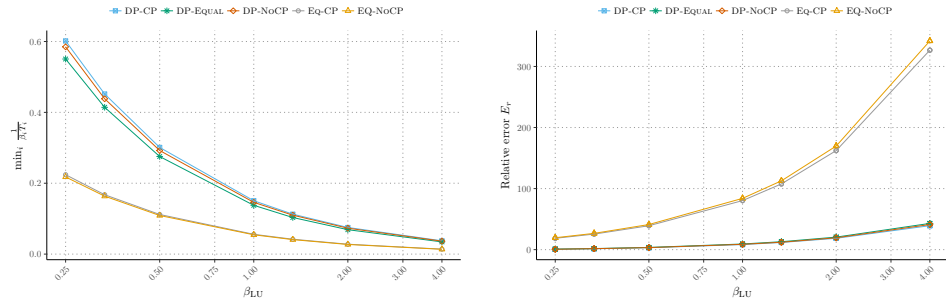


Figure 27: Minimum throughput and relative error for 2MG+LU.

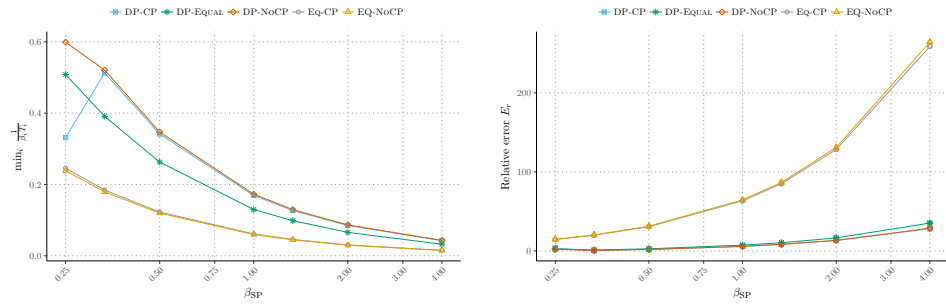


Figure 28: Minimum throughput and relative error for 2MG+SP.

ative error, obviously DP-CP, DP-EQUAL and DP-NoCP perform better than EQ-CP and EQ-NoCP. Among DP-CP, DP-EQUAL and DP-NoCP, we see that the cache partitioning version is the best method to minimize the relative error.

Summary: Overall, we showed that we can obtain important gains using cache partitioning (CP) when co-scheduling three applications, but it is not always the case. The difficulty of obtaining some gain with CP increases with the number of applications involved. The first reason lies in the cache size, often too small to be efficiently partitioned between the applications. The second reason is related to the behavior of the co-scheduled applications. The results show that co-scheduling one or two compute-intensive applications, such as CG, plus one memory-intensive application, such as MG, is a good way to achieve significant improvements with CP. CG is a compute-intensive kernel that performs a lot of irregular memory accesses, while MG is a memory-intensive kernel, hence if we co-schedule one CG and one MG, MG will evict very often cache lines belonging to CG, which will slow down its execution.

8 Conclusion

We have investigated the problem of co-scheduling iterative HPC applications, using the CAT technology provided by Intel to partition the cache. We have proposed a model for the execution time of each application, given a number of cores and a fraction of cache, and we have shown how to instantiate the model on applications coming from the NAS benchmarks. The model turns out to be accurate, as shown in the experiments where we compare the execution time predicted by the model to the real execution time. Several scheduling strategies have been designed, with the goal to maximize the minimum weighted throughput of each application. In particular, we have introduced an optimal strategy for the model, based upon a dynamic programming algorithm. The results demonstrate that in practice, the optimal strategy often leads to better results than a naive strategy sharing equally the resources between applications. Also, we have determined which combinations of applications benefit most from cache partitioning, and demonstrated the usefulness of cache partitioning.

Future work will be devoted to pursuing this experimental study. We hope to get access to platforms with larger shared caches, so that we could scale up the experiments and confirm the usefulness of cache partitioning techniques.

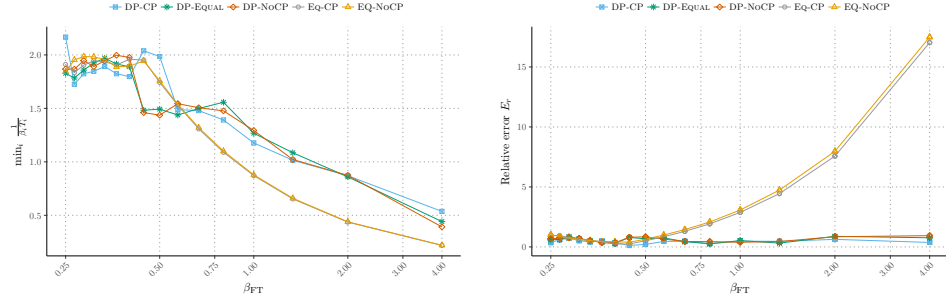


Figure 29: Minimum throughput and relative error for CG, MG and FT.

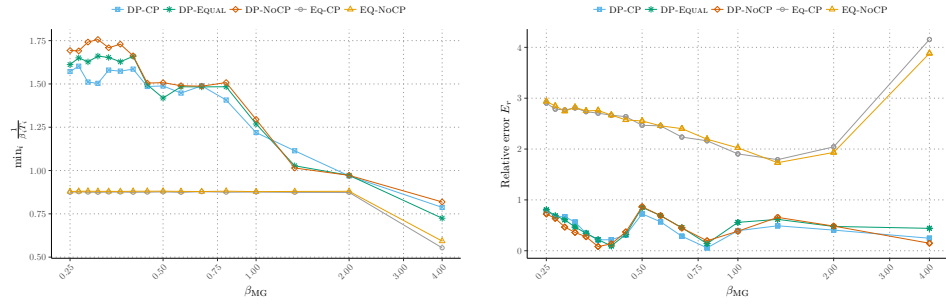


Figure 30: Minimum throughput and relative error for CG, FT and MG.

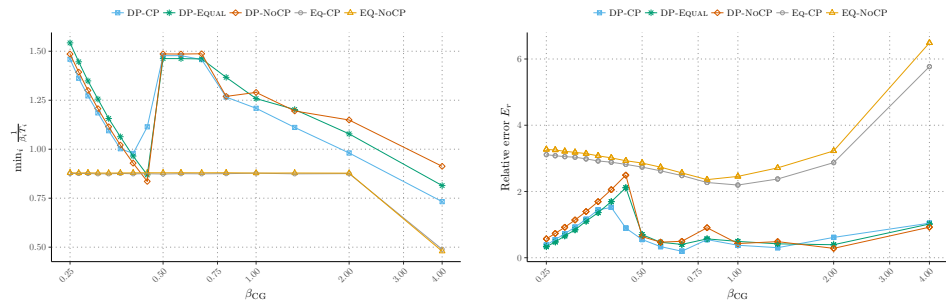


Figure 31: Minimum throughput and relative error for MG, FT and CG.

References

- [1] Muralidhara, S.P., Subramanian, L., Mutlu, O., Kandemir, M., Moscibroda, T.: Reducing memory interference in multicore systems via application-aware memory channel partitioning. In: Proc. 44th IEEE/ACM Int. Sym. Microarchitecture. MICRO-44, ACM (2011) 374–385
- [2] Lo, D., Cheng, L., Govindaraju, R., Ranganathan, P., Kozyrakis, C.: Improving resource efficiency at scale with Heracles. ACM Transactions on Computer Systems (TOCS) **34**(2) (2016)
- [3] Erich Strohmaier et al.: The top500 benchmark (2017) <https://www.top500.org/>.
- [4] Computing, P.: Zettascaler-2.0 configurable liquid immersion cooling system (2017)
- [5] Leverich, J., Kozyrakis, C.: Reconciling high server utilization and sub-millisecond quality-of-service. In: 9th European Conf. on Computer Systems. (2014)
- [6] Zhuravlev, S., Blagodurov, S., Fedorova, A.: Addressing shared resource contention in multicore processors via scheduling. ACM Sigplan Notices **45**(3) (2010) 129–142
- [7] Zhang, Y., Laurenzano, M.A., Mars, J., Tang, L.: Smite: Precise QOS prediction on real-system SMT processors to improve utilization in warehouse scale computers. In: Proc. of the 47th Int. Symp. on Microarchitecture. (2014) 406–418
- [8] Bui, B.D., Caccamo, M., Sha, L., Martinez, J.: Impact of cache partitioning on multi-tasking real time embedded systems. In: 4th IEEE Int. Conf. on Embedded and Real-Time Computing Systems and Applications, IEEE Computer Society (2008) 101–110
- [9] Tian, K., Jiang, Y., Shen, X.: A study on optimally co-scheduling jobs of different lengths on chip multiprocessors. In: Proc. 6th ACM Conf. Computing Frontiers. CF '09, ACM (2009) 41–50
- [10] Nguyen, K.T.: Introduction to Cache Allocation Technology in the Intel® Xeon® Processor E5 v4 Family (February 2016) <https://software.intel.com/en-us/articles/introduction-to-cache-allocation-technology>.
- [11] Sewell, C., et al.: Large-scale compute-intensive analysis via a combined in-situ and co-scheduling workflow approach. In: Proc. of the Int. Conf.

- for High Perf. Computing, Networking, Storage and Analysis, SC'15. (2015)
- [12] Bauer, A.C., Abbasi, H., Ahrens, J., Childs, H., Geveci, B., Klasky, S., Moreland, K., O'Leary, P., Vishwanath, V., Whitlock, B., et al.: In situ methods, infrastructures, and applications on high performance computing platforms. In: *Computer Graphics Forum*. Volume 35., Wiley Online Library (2016) 577–597
- [13] Malakar, P., Vishwanath, V., Munson, T., Knight, C., Hereld, M., Leyffer, S., Papka, M.E.: Optimal scheduling of in-situ analysis for large-scale scientific simulations. In: *Proc. of the Int. Conf. for High Performance Computing, Networking, Storage and Analysis, SC'15*. (2015)
- [14] Dreher, M., Raffin, B.: A Flexible Framework for Asynchronous In Situ and In Transit Analytics for Scientific Simulations. In: *14th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing*, Chicago, United States, IEEE Computer Science Press (May 2014)
- [15] Bao, S., Huo, Y., Parvathaneni, P., Plassard, A.J., Bermudez, C., Yao, Y., Llyu, I., Gokhale, A., Landman, B.A.: A data colocation grid framework for big data medical image processing-backend design. arXiv preprint arXiv:1712.08634 (2017)
- [16] Zhuravlev, S., Saez, J.C., Blagodurov, S., Fedorova, A., Prieto, M.: Survey of scheduling techniques for addressing shared resources in multicore processors. *ACM Computing Surveys (CSUR)* **45**(1) (2012) 4
- [17] Kim, S., Chandra, D., Solihin, Y.: Fair cache sharing and partitioning in a chip multiprocessor architecture. In: *Proceedings of the 13th International Conference on Parallel Architectures and Compilation Techniques*, IEEE Computer Society (2004) 111–122
- [18] Qureshi, M.K., Patt, Y.N.: Utility-based cache partitioning: A low-overhead, high-performance, runtime mechanism to partition shared caches. In: *Microarchitecture, 2006. MICRO-39. 39th Annual IEEE/ACM International Symposium on*, IEEE (2006) 423–432
- [19] Nesbit, K.J., Laudon, J., Smith, J.E.: Virtual private caches. *ACM SIGARCH Computer Architecture News* **35**(2) (2007) 57–68
- [20] Taylor, G., Davies, P., Farmwald, M.: The tlb slice-a low-cost high-speed address translation mechanism. In: *Computer Architecture, 1990. Proceedings., 17th Annual International Symposium on*, IEEE (1990) 355–363

- [21] Tam, D., Azimi, R., Soares, L., Stumm, M.: Managing shared l2 caches on multicore systems in software. In: Workshop on the Interaction between Operating Systems and Computer Architecture, Citeseer (2007) 26–33
- [22] Lin, J., Lu, Q., Ding, X., Zhang, Z., Zhang, X., Sadayappan, P.: Gaining insights into multicore cache partitioning: Bridging the gap between simulation and real systems. In: High Performance Computer Architecture, 2008. HPCA 2008. IEEE 14th International Symposium on, IEEE (2008) 367–378
- [23] Guan, N., Stigge, M., Yi, W., Yu, G.: Cache-aware scheduling and analysis for multicores. In: Proc. 7th ACM Int. Conf. Embedded Software. EMSOFT '09, ACM (2009) 245–254
- [24] Hartstein, A., Srinivasan, V., Puzak, T., Emma, P.: On the nature of cache miss behavior: Is it $\sqrt{2}$. The Journal of Instruction-Level Parallelism **10** (2008) 1–22
- [25] Aupy, G., Benoit, A., Dai, S., Pottier, L., Raghavan, P., Robert, Y., Shantharam, M.: Co-scheduling Amdahl applications on cache-partitioned systems. The Int. Journal of High Performance Computing Applications **32**(1) (2018) 123–138
- [26] Amdahl, G.: The validity of the single processor approach to achieving large scale computing capabilities. In: AFIPS Conference Proceedings. (1967) 483–485
- [27] Krishna, A., Samih, A., Solihin, Y.: Data sharing in multi-threaded applications and its impact on chip design. In: Int. Symp. Performance Analysis of Systems and Software (ISPASS), IEEE (2012) 125–134
- [28] Rogers, B.M., Krishna, A., Bell, G.B., Vu, K., Jiang, X., Solihin, Y.: Scaling the bandwidth wall: challenges in and avenues for CMP scaling. ACM SIGARCH Computer Architecture News **37**(3) (2009) 371–382
- [29] Browne, S., Dongarra, J., Garner, N., Ho, G., Mucci, P.: A portable programming interface for performance evaluation on modern processors. The international journal of high performance computing applications **14**(3) (2000) 189–204
- [30] Bailey, D.H., et al.: The NAS Parallel Benchmarks - Summary and Preliminary Results. In: Proc. of the 1991 ACM/IEEE Conf. on Supercomputing. (1991)



**RESEARCH CENTRE
GRENOBLE – RHÔNE-ALPES**

Inovallée
655 avenue de l'Europe Montbonnot
38334 Saint Ismier Cedex

Publisher
Inria
Domaine de Voluceau - Rocquencourt
BP 105 - 78153 Le Chesnay Cedex
inria.fr

ISSN 0249-6399